# CS5500

## Reinforcement Learning (RL)
### Assignment No. 1

Name - Raktim Gautam Goswami

Roll number - EE17 BTECH11051

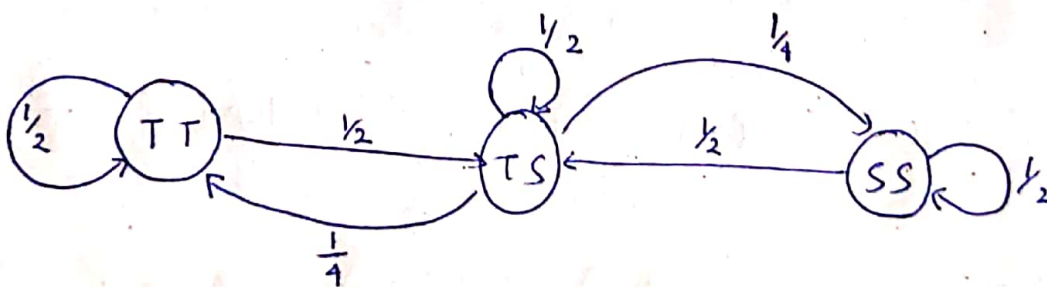**Q1)**

**Sol"-** a) There are 3 states - Tall, short, medium.

~~If initial set of~~

If individual is tall, height of offspring = $\begin{cases} \text{Tall} & \text{w.p. } \frac{1}{2} \\ \text{medium} & \text{w.p. } \frac{1}{2} \\ \text{short} & \text{w.p. } 0 \end{cases}$

If individual is medium, height of offspring = $\begin{cases} \text{Tall} & \text{w.p. } \frac{1}{4} \\ \text{medium} & \text{w.p. } \frac{1}{2} \\ \text{short} & \text{w.p. } \frac{1}{4} \end{cases}$

If individual is short, height of offspring = $\begin{cases} \text{Tall} & \text{w.p. } 0 \\ \text{medium} & \text{w.p. } \frac{1}{2} \\ \text{short} & \text{w.p. } \frac{1}{2} \end{cases}$

∴ the state-space diagram is as follows -



The Transition probability matrix is as follows.

$$P = \begin{array}{c} \\ T \\ M \\ S \end{array} \begin{array}{c} \begin{array}{ccc} T & M & S \end{array} \\ \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \end{array}$$

b) Probabilities that the offspring of

Required probability of first generation offspring

$$= \text{second row of } P \text{ matrix} = \begin{cases} \frac{1}{4} & \text{for tall} \\ \frac{1}{2} & \text{for medium} \\ \frac{1}{4} & \text{for small short} \end{cases}$$

Required probability of second generation offspring

$$= \text{second row of } P^2 \text{ matrix} = \begin{cases} \frac{1}{4} & \text{for tall} \\ \frac{1}{2} & \text{for medium} \\ \frac{1}{4} & \text{for short} \end{cases}$$

Before Similarly,

required probability of third generation offspring

$$= \text{second row of } P^3 \text{ matrix} = \begin{cases} \frac{1}{4} & \text{for tall} \\ \frac{1}{2} & \text{for medium} \\ \frac{1}{4} & \text{for short} \end{cases}$$

c) Clearly, they follows is a trend.

∴ Required probability = second row of $P^n$ matrix

$$= \begin{cases} \frac{1}{4} & \text{for tall} \\ \frac{1}{2} & \text{for medium} \\ \frac{1}{4} & \text{for short.} \end{cases}$$

Q2)

Sol^n. a) The states are

S, 1, 3, 5, 6, 7, 8, W

The Transition matrix (P) is as follows –

|   | S | 1 | 3 | 5 | 6 | 7 | 8 | W |
|---|---|---|---|---|---|---|---|---|
| S | 0 | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | 0 |
| 1 | 0 | 0 | 1/6 | 1/6 | 1/6 | 2/6 | 1/6 | 0 |
| 3 | 0 | 0 | 1/6 | 1/6 | 1/6 | 1/6 | 2/6 | 0 |
| 5 | 0 | 0 | 1/6 | | 1/6 | 1/6 | 1/6 | 1/6 |
| 6 | 0 | 0 | 1/6 | 0 | 2/6 | 1/6 | 1/6 | 1/6 |
| 7 | 0 | 0 | 1/6 | 0 | 0 | 3/6 | 1/6 | 1/6 |
| 8 | 0 | 0 | 1/6 | 0 | 0 | 0 | 4/6 | 1/6 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

b) The only absorbing state is W

c) We know, $V = R + \gamma P V$ ; $\begin{bmatrix} \text{we have } P \text{ from} \\ \text{part (a)} \end{bmatrix}$
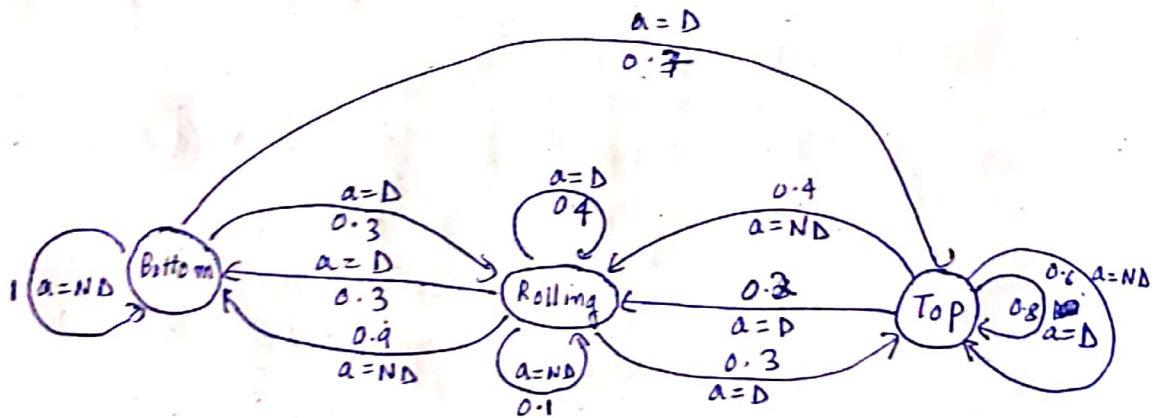
Let $\gamma = 0.9$

$R = [0\ 0\ 0\ 0\ 0\ 0\ 1]^T$, i.e., reward is 0 for all states except for state W.

$$\therefore V = (I - \gamma P)^{-1} R$$

$$= [4.84, 4.94, 4.94, 5.60, 5.60, 5.60, 5.60, 10]^T$$

Q3)

Ans: a)



D → Driving
ND → Not driving

b) A deterministic policy is

$$\pi(s) = \begin{cases} \text{Drive w.p. 1} & \text{when } s = \text{Bottom} \\ \text{Drive w.p. 1} & \text{when } s = \text{Rolling} \\ \text{Don't drive w.p. 1} & \text{when } s = \text{Top} \end{cases}$$
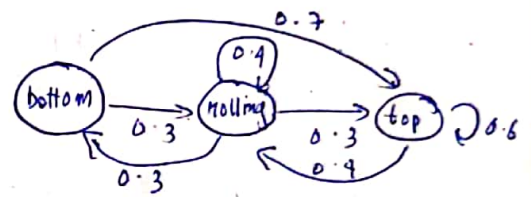
A stochastic policy is

$$\pi(a \mid bottom) = \begin{cases} 0.5 & \text{for } a = \text{Driving} \\ 0.5 & \text{for } a = \text{Not driving} \end{cases}$$

$$\pi(a \mid rolling) = \begin{cases} 0.5 & \text{for } a = \text{Driving} \\ 0.5 & \text{for } a = \text{Not driving} \end{cases}$$

$$\pi(a \mid top) = \begin{cases} 0.5 & \text{for } a = \text{Driving} \\ 0.5 & \text{for } a = \text{Not driving} \end{cases}$$

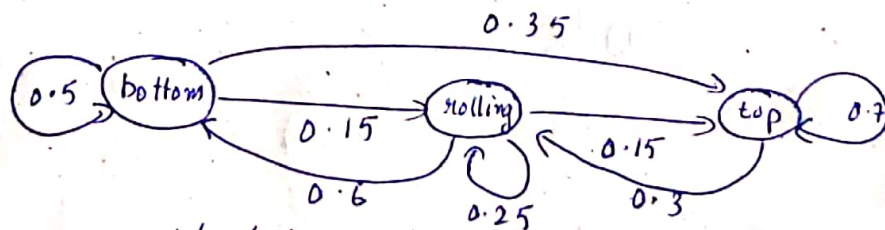c) For deterministic policy, transition probability matrix is

$$P = \begin{array}{c} \\ bottom \\ rolling \\ top \end{array} \begin{array}{c} \overset{bottom \quad rolling \quad top}{} \\ \begin{bmatrix} 0 & 0.3 & 0.7 \\ 0.3 & 0.4 & 0.3 \\ 0 & 0.4 & 0.6 \end{bmatrix} \end{array}$$



For stochastic policy, each element of matrix is determined as follows —

$$P(s' \mid s) = \sum \big( \pi(a \mid s) \, P(s' \mid s, a) \big)$$

$$P = \begin{array}{c} \\ bottom \\ rolling \\ top \end{array} \begin{array}{c} \overset{bottom \quad rolling \quad top}{} \\ \begin{bmatrix} 0.5 & 0.15 & 0.35 \\ 0.6 & 0.25 & 0.15 \\ 0 & 0.3 & 0.7 \end{bmatrix} \end{array}$$



d) A history dependent policy is

$$\pi(a \mid s_t, s_{2}, \dots s)$$

$$\pi(a \mid s_t, s_{t-1}, s_{t-2}, \dots s_1) = \begin{cases} \text{@ } 0.3 & \text{when } H_t < \sum_{i=1}^{t-1} H_i \\ \text{® } 0.4 & \text{when } H_t = \sum_{i=1}^{t-1} H_i \\ \text{© } 0.3 & \text{when } H_t > \sum_{i=1}^{t-1} H_i \end{cases}$$

Q9)

Sol⁻ → for policy $\pi_1$,

$$P_{\pi_1} = \begin{array}{c} \\ A \\ B \\ C \\ D \end{array} \begin{array}{cccc} A & B & C & D \\ \left[\begin{array}{c|c|c|c} 0 & 0.9 & 0.1 & 0 \\ \hline 0.1 & 0 & 0 & 0.9 \\ \hline 0.9 & 0 & 0 & 0.1 \\ \hline 0 & 0 & 0 & 0.1 \end{array}\right] \end{array}$$

For policy $\pi_2$

$$P_{\pi_2} = \begin{array}{c} \\ A \\ B \\ C \\ D \end{array} \begin{array}{cccc} A & B & C & D \\ \left[\begin{array}{c|c|c|c} 0 & 0.1 & 0.9 & 0 \\ \hline 0.9 & 0 & 0 & 0.1 \\ \hline 0.1 & 0 & 0 & 0.9 \\ \hline 0 & 0 & 0 & 1 \end{array}\right] \end{array}$$

For policy $\pi_3$

$$P_{\pi_3} = \begin{array}{c} \\ A \\ B \\ C \\ D \end{array} \begin{array}{cccc} A & B & C & D \\ \left[\begin{array}{c|c|c|c} 0 & 0.42 & 0.58 & 0 \\ \hline 0.1 & 0 & 0 & 0.9 \\ \hline 0.1 & 0 & 0 & 0.9 \\ \hline 0 & 0 & 0 & 1 \end{array}\right] \end{array}$$

$$R = \begin{bmatrix} -10 & -10 & -10 & 100 \end{bmatrix}^T$$

a)   → Assuming $\gamma = 0.9$,

$$V = (1-\gamma P)^{-1} R$$

$$\therefore V^{\pi_1} = \begin{bmatrix} 755 & 867 & 691 & 1000 \end{bmatrix}^T$$

$$V^{\pi_2} = \begin{bmatrix} 755 & 691 & 867 & 1000 \end{bmatrix}^T$$

$$V^{\pi_3} = \begin{bmatrix} 772 & 869 & 869 & 1000 \end{bmatrix}^T$$

b) $\pi_3$ seems to be the best one as they all the values of $V^{\pi_3}$ are higher compared to others

Scanned by CamScanner

(top right margin)
0.9×0.9
+0.6×0.1

0.4×0.1
+0.1×0.1

c) $\pi_2$ & $\pi_1$ are not ~~of it~~ comparable as here in $V^{\pi_2}$ & $V^{\pi_3}$ two ~~~~ states are having same values & in other 2 states, the values are interchanged. However, $\pi_3$ can be compared as $V^{\pi_3}$ has all elements greater than $V^{\pi_2}$ & $V^{\pi_1}$ ∴ it is the best policy.

Q5)

$sol^n$ a) We can do this using either value iteration or policy iteration. In this case, we will solve using value iteration.

Step I:

We initialise $V_1(s)$ for $s \in S$, @ a small value $\epsilon$.

Step II:

For all states, we find

$$V_{k+1}(s) \longleftarrow \max_a \left[ \sum_{s' \in S} P_{ss'}^a \left( R_{ss'}^a + \gamma V_k(s') \right) \right]$$

Step III:

If $|V_{k+1}(s) - V_k(s)| < \epsilon$ for all $s \in S$, we go to step IV, $\quad$ (Let $V_{k+1}(s) = V_*(s)$)

else,

we go back to step II.

Step IV:

For all states $s \in S$,

$$\pi_*(s) = \underset{a}{\arg\max} \; V_*(s)$$

∴ We get $\pi_*(s)$ such that $V^{\pi_*}(s)$ is maximum for all states of the MDP.

Q6) The python code for this question is attached with the submission

sol:- a) For $\gamma = 1$, we cannot find an optimal value or policy as it becomes an infinite process and value function becomes unbounded

b) For $\gamma = 0.9$,

$$V_* = [65.6, 72.8, 80.9, 89.9, 99.9, 99.9]^T$$

using $\pi^* = \max\limits_{a} V_*^a$

$$\pi^*(s_1) = \begin{cases} 1 & \text{for moving right} \\ 0 & \text{" " left} \end{cases}$$

$$\pi^*(s_2) = \pi^*(s_3) = \pi^*(s_4) = \pi^*(s_5) = \begin{cases} 1 & \text{for moving right} \\ 0 & \text{" " left} \end{cases}$$

$$\pi^*(s_6) = 1 \quad \text{for staying at } s_6$$

For $\gamma = 0.5$,

$$V_* = \begin{bmatrix} 1.24 & 2.49 & 4.99 & 9.99 & 19.99 & 19.99 \end{bmatrix}^T$$

$\pi^*(s)$ is same as for $\gamma = 0.9$

for $\gamma = 0.1$,

$$V_* = \begin{bmatrix} 0.0011 & 0.011 & 0.11 & 1.1 & 11 & 11 \end{bmatrix}^T$$

$\pi^*(s)$ remains same as for $\gamma = 0.9$

∴ it can be observed that the optimal value function changes for different $\gamma$. However, policy remains the same.

c) Adding $c$ to all rewards is found out by running the program in python.

For $c = -1$, $\gamma = 0.9$,

$$V_* = \begin{bmatrix} 55.6 & 62.8 & 70.9 & 79.9 & 89.9 & 89.9 \end{bmatrix}^T$$

for $c = 1$, $\gamma = 0.9$,

$$V_* = \begin{bmatrix} 75.6 & 82.8 & 90.9 & 99.9 & 109.9 & 109.9 \end{bmatrix}^T$$

for $c = 10$, $\gamma = 0.9$

$$V_* = \begin{bmatrix} \cancel{75.6} & \cancel{82.8} & \cancel{90.9} & \cancel{99.9} & \cancel{109.9} & \cancel{109} \end{bmatrix}$$

$$V_* = \begin{bmatrix} 165.6 & 172.8 & 180.9 & 189.9 & 199.9 & 199.9 \end{bmatrix}^T$$

∴ the same value reward is being added to all the rewards, the policy remains the same.

d) for any policy $\pi$,

$$V^\pi = (I - \gamma P)^{-1} R$$

$$\hat{v}^{\pi} = (1-\gamma P)^{-1}(R+C) \quad ; \quad \text{when } C = [\underbrace{c \ c \ \ldots \ c}_{\text{No. of state}}]^T$$

$\therefore$ clearly, $\hat{v}^{\pi} = v^{\pi} + (1-\gamma P)^{-1} C$

08)

Sol$^n$:

$$L(v) = \max_{a \in A} \left[ R^a + \gamma P^a v \right]$$

$\therefore$ $V_*$ is a fixed point of operator $L$,

$\therefore$ $T(v_*^\rho) = V_*$

$$|V_{k+1} - V_*|_\infty = |T(V_k) - T(V_*)|_\infty$$

$$= \left| \max_{a \in A} [R^a + \gamma P^a V_k] - \max_{a \in A} [R^a + \gamma P^a V_*] \right|_\infty$$

$$\leq \max_{a \in A} \left| [R^a + \gamma P^a V_k] - [R^a + \gamma P^a V_*] \right|_\infty$$

$$= \gamma \left| P^a (V_k - V_*) \right|_\infty$$

$$\leq \gamma \left| V_k - V_* \right|_\infty$$

$\therefore$ $\left| V_{k+1} - V_* \right|_\infty \leq \gamma \left| V_k - V_* \right|_\infty$

$$\leq \gamma^2 \left| V_{k-1} - V_* \right|_\infty$$

$$\leq \gamma^3 \left| V_{k-2} - V_* \right|_\infty$$

$$\leq \gamma^k \left| V_1 - V_* \right|_\infty$$

$\therefore$ $\left| V_{k+1} - V_* \right|_\infty \leq \gamma^k \left| V_1 - V_* \right|_\infty$

7)

Sol<sup>n</sup>—

The value for any state can be found out by recursively running

$$V^{\pi}(s) \leftarrow \max_{a \in A} \left\{ \sum_{s' \in S} P_{ss'}^a \left[ R_{ss'}^a + \gamma V^{\pi}(s') \right] \right\}$$

This is carried on for all the states with different values of $\gamma$ and $P$. $(\because$ on changing noise, $P$ changes$)$

$\therefore$ The required parameters are as follows —

For close exit, risking the cliff,        $\gamma = 0.1$, noise $= 0.45$

For distant exit, risking the cliff,      $\gamma = 0.99$, noise $= 0.5$

For close exit, avoiding the cliff,       $\gamma = 0.99$, noise $\approx 0$

For distant exit, avoiding the cliff,     $\gamma = 0.1$, noise $= 0$